

## A Few Remarks on the Index of Context-free Grammars and Languages

J. GRUSKA\*

*Department of Computer, Information, and Control Sciences  
University of Minnesota*

A hierarchy of context-free grammars and languages with respect to the index of context-free grammars is established and the undecidability of the basic problems is proven.

### 1. NOTATION

Let  $G = \langle V, \Sigma, P, \sigma \rangle$  be a context-free grammar (in short, a grammar) where  $\Sigma \subset V$  is the set of terminals,  $V - \Sigma$  the set of nonterminals  $\sigma \in V - \Sigma$  is the initial symbol of  $G$  and  $P \subset (V - \Sigma) \times V^*$  is a finite set of productions of  $G$ .  $L(G) = \{x; \sigma \xrightarrow{*} x \in \Sigma^*\}$  is the context-free language (in short, a language) generated by  $G$ .

Let  $\epsilon$  denote the empty word,  $|x|$  the length of word  $x$ , and  $I$  the set of positive integers.

Following Brainerd (1968), the index of a derivation  $\tau$ , in short  $\text{Ind}(\tau)$ , where

$$\tau : w_1, w_2, \dots, w_k,$$

is the smallest integer  $i_0$  such that neither of the words  $w_i$ ,  $1 \leq i \leq k$ , has more than  $i_0$  occurrences of nonterminals. For an  $x \in L(G)$ ,  $\text{Ind}(x) = \min\{\text{Ind}(\tau); \tau \text{ is a derivation of } x \text{ from } \sigma \text{ in } G\}$ .

For a grammar  $G$  and a language  $L$  let  $\text{Ind}(G) = \max\{\text{Ind}(x); x \in L(G)\}$ ;  $\text{Ind}(L) = \min\{\text{Ind}(G); L(G) = L\}$ . If  $\text{Ind}(G) < \infty$  ( $\text{Ind}(L) < \infty$ ), then the grammar  $G$  (the language  $L$ ) is said to be of finite index; otherwise of infinite index.

\* Present Address: Mathematical Institute, Slovak Academy of Sciences, Bratislava, Czechoslovakia.

## 2. SALOMAA'S PROBLEM

In a recent paper, Salomaa (1969) raised the question of whether there are two grammars which generate the same language but only one of them is of finite index. The answer is in the positive. Indeed, by Salomaa (1969), there exists a grammar  $G_1$  with the infinite index which generates the Dyck language  $L_0$  over the alphabet  $\{0, 1\}$ .  $L_0$  is a deterministic language and therefore a grammar  $G_2$  generating  $\{0, 1\}^* - L_0$  exists. Combining these two grammars we get a grammar with the infinite index which generates the language  $\{0, 1\}^*$  having the index 1.

## 3. FINITE INDEX LANGUAGES

By Salomaa (1969), there is a context-free language of infinite index. Languages of finite index form a very natural class of languages which has been studied in several papers under different names (superlinear languages, derivation-bounded languages, semilinear languages and so on).

**THEOREM 1.** *The class of finite index languages is the small full AFL which contains linear languages and is closed under substitution.*

*Proof.* Each finite index language is a derivation-bounded language and by Ginsburg and Spanier (1968) the class of derivation-bounded languages forms a full AFL mentioned in the Theorem. On the other hand, each derivation-bounded language can be obtained from linear languages by substitution and, therefore, is of a finite index.

Several infinite hierarchies of finite index languages have recently appeared in the literature (Greibach, 1969; Gruska, 1969). In the following theorem, a new hierarchy depending on the index of languages is proven.

**THEOREM 2.** *For every  $n \in I \cup \{\infty\}$  there is a language  $L_n$  such that  $\text{Ind}(L_n) = n$ .*

*Proof.* Let  $L_0$  be the Dyck language over the alphabet  $\{0, 1\}$ ; i.e., the language generated by the grammar

$$\sigma \rightarrow 0\sigma 1, \quad \sigma \rightarrow \sigma\sigma, \quad \sigma \rightarrow \epsilon.$$

By Salomaa (1969),  $\text{Ind}(L_0) = \infty$ . Trivially,  $\text{Ind}(\{\epsilon\}) = 1$ . For  $n$  finite and  $n > 1$ , let  $L_n = L_0 \cap (0^*1^*)^{2^{n-1}}$ . The theorem will be proved by showing that

$\text{Ind}(L_n) = n$ . Since  $L_n$  is generated by the grammar  $G_n$  with the initial symbol  $\sigma_1$  and the rules<sup>1</sup>

$$\begin{aligned} \sigma_i &\rightarrow 0\sigma_i1 \mid \sigma_{i+1}\sigma_{i+1} \mid \epsilon, & 1 \leq i \leq n-1, \\ \sigma_n &\rightarrow 0\sigma_n1 \mid \epsilon, \end{aligned}$$

we have immediately  $\text{Ind}(L_n) \leq n$ . To complete the proof, it remains to show that  $\text{Ind}(L_n) \geq n$ .

To that end, let  $G = \langle V, \{0, 1\}, P, \sigma \rangle$  be a grammar such that  $\text{Ind}(G) = \text{Ind}(L_n)$ ,  $L(G) = L_n$  and, moreover, all rules of  $G$  have either the form  $A \rightarrow uBv$  or  $A \rightarrow uBCv$  with  $u, v$  being terminal words and  $B, C$  being nonterminals. Clearly, such a  $G$  does exist. We can also assume that  $G$  is a reduced grammar, i.e., all nonterminals are reachable and generate some terminal words.

Let  $m = \max\{|\alpha|; A \rightarrow \alpha \text{ is in } P\}$  and let  $n_0$  be the number of nonterminals of  $G$ .

Let  $\psi : \{0, 1\}^* \rightarrow \{0, 1\}^*$  be the mapping defined by  $\psi(x) = \psi(y_1y_2)$ , if  $x = y_101y_2$  and  $\psi(x) = x$  otherwise.

If  $x \in L_n$  and  $x = yz$ , then  $\psi(y) \in \{0\}^*$  and  $\psi(z) \in \{1\}^*$ . From that and from the structure of words in  $L_n$ , it follows

$$\text{if, in } G, A \xRightarrow{*} xAy \text{ for a nonterminal } A \text{ and terminal words } x, y, \text{ then} \quad (1)$$

$$\psi(x) = 0^k, \psi(y) = 1^k \text{ for some } k \geq 0.$$

If in a derivation tree of a word  $x$  no path contains a nonterminal twice, then  $|x| \leq m^{n_0}$ . In view of (1) this in turn implies

$$\text{if } A \xRightarrow{*} x \text{ in } G, \quad \text{then } |\psi(x)| \leq m^{n_0}, \quad (2)$$

and, as a corollary,

$$\text{if } A \xRightarrow{*} x \in \{1\}^*\{0\}^*, \quad \text{then } |x| \leq m^{n_0}. \quad (3)$$

Now let  $N > 3m^{n_0+1}$  be an integer and let  $y_i, i \geq 0$ , be words defined by

$$y_0 = \epsilon, \quad y_{i+1} = 0^{2N}y_i1^N0^Ny_i1^{2N}, \quad i \geq 0.$$

For every  $A$  in  $V$ , let  $\nu(A) = \{A\} \cup \{x; A \xRightarrow{*} x \in \{1\}^*\{0\}^*\}$ . By (3),  $\nu(A)$  is a finite set.

<sup>1</sup> The productions are written in an abbreviated form  $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$  instead of  $A \rightarrow \alpha_1, A \rightarrow \alpha_2, \dots, A \rightarrow \alpha_n$ .

Let us now construct a new grammar  $G'$  from  $G$  by replacing every production  $A \rightarrow \alpha$  of  $G$  by the set of productions  $\{A \rightarrow z; z \in \nu(\alpha)\}$ , where  $\nu(xa) = \nu(x)\nu(a)$ , for any word  $x$  and symbol  $a$ .

Clearly,  $L(G') = L_n$ ,  $\text{Ind}(G') = \text{Ind}(L_n)$ . The length of right sides of productions of  $G'$  is not more than  $m^{n_0+1}$  and (1)–(3) hold also for  $G'$ . For the rest of this proof we will deal only with the grammar  $G'$  and, therefore, all derivation concepts, for example  $\xrightarrow{*}$ , refer to  $G'$ .

Assume now for a moment that the following lemma has already been proved.

**LEMMA 3.** *Let  $A \xrightarrow{*} xy_i z$ , where  $A$  is a nonterminal,  $i \geq 0$  and either  $x = 0^v$ ,  $z = 1^\mu$  or  $x = 0^v$ ,  $z = 1^N 0^\mu$  or  $x = 1^v 0^N$ ,  $z = 1^\mu$ , with  $v, \mu \leq N$  being integers. Then  $\text{Ind}(\tau) \geq i + 1$  for any derivation  $\tau$  of  $xy_i z$  from  $A$ .*

This yields immediately  $\text{Ind}(G') \geq n$  and thus  $\text{Ind}(L_n) \geq n$ .

Hence, to complete the proof of the Theorem, it remains only to prove the Lemma. The proof will be by induction on  $i$ . The case  $i = 0$  is trivial. Assume that Lemma holds for  $0, 1, \dots, i - 1$ , and let

$$A = W_0, W_1, \dots, W_s = xy_i z \tag{4}$$

be a derivation of  $xy_i z$  from  $A$  of the minimal index and as short as possible.

Since  $N > 3m^{n_0+1}$ , (2) implies that there exists the smallest integer  $i_0$  such that  $W_{i_0+1}$  contains two nonterminals. Thus,  $W_{i_0} = u_{i_0} A_{i_0} v_{i_0}$ , where  $u_{i_0} v_{i_0} \in \{0, 1\}^*$  and  $A_{i_0}$  is a nonterminal.

We claim that

$$|u_{i_0}| \leq |x| + N + m^{n_0+1}, \quad |v_{i_0}| \leq |z| + N + m^{n_0+1}. \tag{5}$$

Assume that (5) does not hold. Then there must exist the smallest integer  $i_1 \leq i_0$  such that  $W_{i_1} = u_{i_1} A_{i_1} v_{i_1}$ ,  $u_{i_1} v_{i_1} \in \{0, 1\}^*$  and

$$\text{either } |u_{i_1}| > |x| + N + m^{n_0+1} \quad \text{or } |v_{i_1}| > |z| + N + m^{n_0+1}. \tag{6}$$

Let  $z_{i_1}$  be such that  $u_{i_1} z_{i_1} v_{i_1} = xy_i z$ . Since  $y_i = 0^{2N} y_{i-1} 1^N 0^N y_{i-1} 1^{2N}$  and  $N > 3m^{n_0+1}$ , from (6) it follows that  $\psi(z_{i_1}) > m^{n_0}$  what contradicts to (2). Thus (5) holds.

Using (5), we can now complete the proof of Lemma. Clearly,  $W_{i_0+1} = u_{i_0} u' B C v' v_{i_0}$  for some terminal words  $u'$  and  $v'$ . Since (4) was the shortest derivation of  $z$  from  $A$  among those with minimal index, none of the words  $\bar{u}, \bar{v}$  where  $B \xrightarrow{*} \bar{u}$ ,  $C \xrightarrow{*} \bar{v}$ ,  $z = u_{i_0} u' \bar{u} \bar{v} v' v_{i_0}$  is in  $\{1\}^* \{0\}^*$ . But it means that

$\bar{u}$  and  $\bar{v}$  are of the form  $x_1 y_{i-1} z_1$  with either  $x_1 = 0^{\nu_1}$ ,  $z_1 = 1^{\mu_1}$  or  $x_1 = 0^{\mu_1}$ ,  $z_1 = 1^{\nu_1} 0^{\mu_1}$  or  $x_1 = 1^{\nu_1} 0^{\mu_1}$ ,  $z_1 = 1^{\mu_1}$  and  $\nu_1, \mu_1 \leq N$ . By induction hypothesis  $\text{Ind}(\tau_1) \geq i$ ,  $\text{Ind}(\tau_2) \geq i$  for any derivations  $\tau_1$  of  $\bar{u}$  from  $B$  and  $\tau_2$  of  $\bar{v}$  from  $C$ . Hence  $\text{Ind}(\tau) \geq i + 1$ , completing the proof.

4. UNDECIDABILITY

For any  $n \in I \cup \{\infty\}$  there exist a context-free grammar  $G_n$  and a context-free language  $L_n$  such that  $\text{Ind}(G_n) = n = \text{Ind}(L_n)$ . For  $n$  infinite it follows from Salomaa's result (1969). For  $n$  finite, the existence of  $L_n$  was proved in the previous section and as  $G_n$  we can take the grammar with the rules  $\sigma \rightarrow A^n, A \rightarrow a$ .

On the other hand, as it will be shown in this section, for any  $k \in I \cup \{\infty\}$  it is undecidable for a context-free grammar  $G$  whether or not  $\text{Ind}(G) = k$  or whether or not  $\text{Ind}(L(G)) = k$ .

Two more results are proved in this section. If we think of  $\text{Ind}$  as being a criterion of complexity of grammars and languages, then they may be interpreted as follows: (i) It is undecidable whether a given grammar  $G$  is a simplest grammar for  $L(G)$ ; (ii) There is no effective way to find a simplest grammar for  $L(G)$ , given a grammar  $G$ .

All these results will follow easily from the following lemma.

To simplify the ensuing discussion, let us denote by  $P(x, y)$  the predicate which holds true if and only if  $x$  and  $y$  are  $n$ -tuples of nonempty words for some integer  $n$  and the post-correspondence problem for  $x$  and  $y$  has a solution.

LEMMA 4. For any  $n$ -tuples  $x = (x_1, \dots, x_n), y = (y_1, \dots, y_n)$  of nonempty words over  $\Sigma = \{0, 1\}$ , and context-free grammars  $G_1$  and  $G_2$  with  $L_i = L(G_i)$ , a grammar  $G$  can be effectively found such that  $L(G) = L' \cup L''$ , where

$$L' = \{ucwcv^R; u, v \in \Sigma, u \neq v, w \in L_1\}$$

and

$$L'' = \{x_{i_1} \dots x_{i_k} c w c y_{i_k}^R \dots y_{i_1}^R; k \geq 1, w \in L_2\}^2.$$

Furthermore, if  $L_2 \subseteq L_1$ , then

$$\begin{aligned} \text{Ind}(G) &\leq \text{Ind}(G_1) \text{ if } P(x, y) \text{ does not hold} \\ &= \max\{\text{Ind}(G_1), \text{Ind}(G_2)\} \text{ if } P(x, y) \text{ holds,} \end{aligned}$$

<sup>2</sup> For a word  $x = x_1 x_2 \dots x_n, x^R = x_n \dots x_2 x_1$ .

and if, moreover,  $\text{Ind}(L_1) \leq \text{Ind}(L_2) = \text{Ind}(G_2)$ , then

$$\begin{aligned} \text{Ind}(L(G)) &\leq \text{Ind}(L_1) \text{ if } P(x, y) \text{ does not hold} \\ &= \text{Ind}(L_2) \text{ if } P(x, y) \text{ holds.} \end{aligned}$$

*Proof.* Let  $G_i = \langle V_i, \Sigma, P_i, S_i \rangle$ ,  $V_1 \cap V_2 = \emptyset$ ,  $A_1, A_2, A_3, S', S''$ ,  $S \notin V_1 \cup V_2$ .

Let  $G$  have the initial symbol  $S$  and productions  $P_1 \cup P_2$  together with

$$\begin{aligned} S &\rightarrow S' \mid S'', \\ S' &\rightarrow 0S'0 \mid 1S'1 \mid 0A_1 \mid 1A_1 \mid 0A_21 \mid 1A_20 \mid A_30 \mid A_31, \\ A_1 &\rightarrow 0A_1 \mid 1A_1 \mid cS_1c, \\ A_2 &\rightarrow A_20 \mid A_21 \mid A_1, \\ A_3 &\rightarrow A_30 \mid A_31 \mid cS_1c, \\ S'' &\rightarrow x_i S'' y_i^R \mid x_i c S_2 c y_i^R, \quad 1 \leq i \leq n. \end{aligned}$$

Clearly,  $S'$  generates  $L'$  and  $S''$  generates  $L''$ .

Let now  $L_2 \subseteq L_1$ . If  $P(x, y)$  does not hold, then  $L'' \subseteq L'$  and, therefore,  $\text{Ind}(G) \leq \text{Ind}(G_1)$ . If  $P(x, y)$  holds, then  $\text{Ind}(G) = \max\{\text{Ind}(G_1), \text{Ind}(G_2)\}$ .

To prove the last assertion of the Lemma, we will use the fact that  $\text{Ind}(L \cap R) \leq \text{Ind}(L)$  if  $L$  is a CFL and  $R$  a regular set. It can be shown easily going through a standard proof of the theorem that the intersection of a language  $L$  and a regular set  $R$  is again a language (see, for example, Ginsburg, 1966).

Let now  $\text{Ind}(L_1) \leq \text{Ind}(L_2) = \text{Ind}(G_2)$ . If  $P(x, y)$  does not hold, then  $\text{Ind}(L(G)) \leq \text{Ind}(L_1)$ .

If  $P(x, y)$  holds, then there are indices  $i_1, \dots, i_k$  such that

$$x_{i_1} x_{i_2} \cdots x_{i_k} = y_{i_1} y_{i_2} \cdots y_{i_k}.$$

Consider now the regular set

$$R = x_{i_1} x_{i_2} \cdots x_{i_k} \{0, 1\}^* c y_{i_k}^R \cdots y_{i_2}^R y_{i_1}^R.$$

The intersection of  $R$  and  $L(G)$  has the form

$$L(G) \cap R = x_{i_1} \cdots x_{i_k} c L_2 c y_{i_k}^R \cdots y_{i_1}^R.$$

One can easily prove that if  $L$  is a language and  $a$  a symbol, then  $\text{Ind}(L) = \text{Ind}(aL) = \text{Ind}(La)$ . Thus  $\text{Ind}(L(G) \cap R) = \text{Ind}(L_2)$  and we have  $\text{Ind}(L(G)) \geq \text{Ind}(L_2)$ . On the other hand, from the construction of  $G$  it follows immediately

that  $\text{Ind}(L(G)) \leq \text{Ind}(G) \leq \text{Ind}(G_2) = \text{Ind}(L_2)$ . Thus  $\text{Ind}(L(G)) = \text{Ind}(L_2)$  and this completes the proof of the Lemma.

**THEOREM 5.** *Let  $n \in I \cup \{\infty\}$ . It is undecidable for an arbitrary grammar  $G$  whether or not  $\text{Ind}(G) = n$ .*

**THEOREM 6.** *Let  $n \in I \cup \{\infty\}$ . It is undecidable for an arbitrary grammar  $G$  whether or not  $\text{Ind}(L(G)) = n$ .*

For  $n > 1$ , the theorems follow from the Lemma 4 by taking  $G_1$  to be the grammar  $S_1 \rightarrow S_1 0$ ,  $S_1 \rightarrow S_1 1$ ,  $S_1 \rightarrow \epsilon$  and  $G_2$  to be a grammar with  $\text{Ind}(G_2) = \text{Ind}(L(G_2)) = n$ . As a byproduct we get the theorems for  $n = 1$ .

**COROLLARY 7.** *There is no effective way to determine  $\text{Ind}(L(G))$  for an arbitrary grammar  $G$ .*

**THEOREM 8.** *It is undecidable for an arbitrary grammar  $G$  whether or not  $\text{Ind}(G) = \text{Ind}(L(G))$ .*

*Proof.* Let us take as  $G_1$  the grammar with the rules  $S_1 \rightarrow S_1 S_1$ ,  $S_1 \rightarrow 0$ ,  $S_1 \rightarrow 1$ ,  $S_1 \rightarrow \epsilon$  and as  $G_2$  a grammar such that  $\text{Ind}(G_2) = \text{Ind}(L(G_2)) = \infty$ . Now  $\text{Ind}(G) = \text{Ind}(L(G))$  for the grammar  $G$  from the Lemma 4 if and only if  $P(x, y)$  holds. Hence the Theorem.

By using the same construction as in the proof of foregoing theorem, we get

**COROLLARY 9.** *There is no effective way to construct for an arbitrary grammar  $G$ , a grammar  $G'$  such that  $L(G) = L(G')$  and  $\text{Ind}(G') = \text{Ind}(L(G))$ .*

## 5. MODIFICATION

The index of a grammar  $G$  represents the maximal number of nonterminals which may occur simultaneously in derivation steps in the derivations of elements in  $L(G)$ . However, there is no restriction as to how the nonterminals are spread out in words. In this section, we shall try to put some restriction on the distance between two nonterminals in a derivation step.

By  $\text{Ind}'(G)$  we will mean the smallest integer  $k$  (if such a  $k$  does exist; otherwise we put  $\text{Ind}'(G) = \infty$ ) such that for every  $x \in L(G)$  there is a derivation  $\sigma = w_0, w_1, \dots, w_k = x$  in  $G$  such that each  $w_i = u_i \alpha_i v_i$ , where  $u_i v_i$  is a terminal word and  $|\alpha_i| \leq k$ . Let  $\text{Ind}'(L) = \min\{\text{Ind}'(G); L(G) = L\}$  for a language  $L$ .

Clearly, for every  $n \in I \cup \{\infty\}$  there is a grammar  $G_n$  such that  $\text{Ind}'(G_n) = n$ . Using the technique of the proof of Theorem 5, one can show that if  $n \in I \cup \{\infty\}$ , then it is undecidable for an arbitrary grammar  $G$  whether  $\text{Ind}'(G) = n$ . Since  $\text{Ind}(L) \leq \text{Ind}'(L)$  for every language  $L$ ,  $\text{Ind}'(L_0) = \infty$ . However, for every language  $L$  either  $\text{Ind}'(L) = \infty$  or  $\text{Ind}'(L) = 1$ , and, therefore, the criterion  $\text{Ind}'$  does not induce an infinite hierarchy of context-free languages. In order to show that  $\text{Ind}'(L) < \infty$  implies  $\text{Ind}'(L) = 1$  one can proceed as follows. Let  $G$  be a grammar such that  $\text{Ind}'(G) = k < \infty$  and  $G = \langle V, \Sigma, P, \sigma \rangle$ . Let us form a new grammar  $G' = \langle V', \Sigma, P', \sigma \rangle$  by taking as new nonterminals the symbols  $[\alpha]$ , where  $\alpha \in V^*$ ,  $|\alpha| \leq k$  and the first and the last symbol of  $\alpha$  are nonterminals. If, in  $G$ ,  $\alpha \Rightarrow uv\beta$ ,  $|\beta| \leq k$ ,  $uv \in \Sigma^*$ ,  $\beta$  starts and ends with nonterminal symbols, then we put to  $P'$  the production  $[\alpha] \rightarrow u[\beta]v$  and all productions of  $P'$  are formed in that way. Clearly,  $L(G') = L(G)$ ,  $\text{Ind}'(G') = 1$ . Therefore,  $\text{Ind}'(L) < \infty$  if and only if  $L$  is a linear language. From that and from Greibach (1966) it follows that it is undecidable for an arbitrary grammar  $G$  whether  $\text{Ind}'(L(G)) = \infty$  (or  $\text{Ind}'(L(G)) = 1$ ).

*Added in proof:* Salomaa's problem was solved also by N. D. Jones in *Information and Control* 16, 201–202.

#### ACKNOWLEDGMENT

The author wish to thank the referee for his suggestions concerning the exposition of Section 4.

RECEIVED: June 15, 1970; REVISED: March 8, 1971

#### REFERENCES

- BRAINERD, B. (1968), An analog of a theorem about context-free languages, *Information and Control* 11, 561–567.
- GINSBURG, S. (1966), The mathematical theory of context-free languages, McGraw-Hill, New York.
- GINSBURG, S., AND SPANIER, E. H. (1968), Derivation-bounded Languages, *J. Comput. System Sci.* 2, 228–250.
- GREIBACH, S. A. (1966), The unsolvability of the recognition of linear context-free languages, *J. Ass. Comput. Mach.* 12, 42–52.
- GREIBACH, S. A. (1969), An infinite hierarchy of context-free languages, *JACM* 16, 91–106.
- GRUSKA, J. (1969), Some classifications of context-free languages, *Information and Control* 14, 152–173.
- SALOMAA, A. (1969), On the index of context-free grammars and languages, *Information and Control* 14, 474–477.